

NLP 감정분석 기반
마케팅 시장 분석

키워드에 따른 커뮤니티 별 온라인 여론 감정 분석 시스템

빅데이터미네이터

“

키워드에 따른 커뮤니티 별 온라인 여론 감정 분석 시스템



사용자가 입력하는 키워드 기반으로 커뮤니티의 본문 댓글에 대한 크롤링을 진행하고,

이에 따른 다양한 의견들을 감성 분석하여 결과를 도출하고자 합니다.

이는 온라인 여론에 대한 감정 분석을 진행하는 것으로 마케팅 측면에서도 도움이 되는 자료가 될 것입니다.

1. Input

1. 사용자 input
: 키워드

2. 크롤링 범위
: 제목, 본문, 날짜, 댓글

2. 커뮤니티 선정

- 네이트판
- 뽀뿌
- 클리앙
- 루리웹
(+ 네이버 블로그 및 sns 고려)

: 커뮤니티의 순위, 사용자 수,
카테고리의 다양성, 데이터의 양을
중심으로 선정

3. 전처리

- 영어, 특수문자 등 제거
- 초성 (ㅋㅋ, ㅎㅎ, 초성 욕설 등) 제거
- 형태소 단위로 토큰화 → 문장 단위로 순서는 유지해야함
(앞뒤 단어와의 연관성에 따른 감성 분석 가능)
- csv 파일로 저장 (프로토타입)

4. 프로토타입 구조

크롤링부 (+ 전처리)

↓

csv 파일 (추후 DB 구축)

↓

감성어 사전
(범용 사전 + 감성어 추가)

감성 분석기 => 결과 도출

새롭게 정한 전처리 기준에 맞게 각 커뮤니티에 대한 크롤링부를 구현하고자 합니다.

따라서 각 커뮤니티 별로 크롤러를 구현하는 작업을 시작하였습니다.

커뮤니티 별 크롤링

이. 네이트판



+

```
In [21]: import os
import sys
import pandas as pd
import requests
import re
from datetime import datetime
from bs4 import BeautifulSoup
from selenium.webdriver import Chrome
from selenium.webdriver.common.keys import Keys
import time
import json
import csv

In [22]: browser = Chrome()
browser.maximize_window()
base_url = 'https://nann.nate.com'
browser.get(base_url)
browser.find_elements_by_xpath('//*[@id="input_search"]')[0].click()

In [23]: #원하는 검색어 입력
query_txt = input('크롤링할 키워드를 입력하세요: ')

크롤링할 키워드를 입력하세요: 코로나

In [24]: browser.find_elements_by_xpath('//*[@id="input_search"]')[0].send_keys(query_txt)

In [25]: #검색
browser.find_elements_by_xpath('//*[@id="search"]/fieldset/button')[0].click()

In [26]: #제목, 게시물, 검색 결과, 더 보기 선택
browser.find_elements_by_xpath('//*[@id="container"]/div[2]/div[1]/div[3]/p/a')[0].click()
#우주 최신순 선택 추가

In [35]: nate_page = browser.current_url
page = nate_page + '&page='
#page 주소 부분 keyword 입력 후 주소 바꿈은 글씨로 바꾸기
titles = []
links = []
#page 1-4까지 추출할분 (10page까지 크롤링시 시간도 다소 길리럼, 본문의 'p'와 'espresso_editor_view'등의 class에 대한 부분 추가해야)
for i in range(1, 6):
    res = requests.get(page + str(i))
    res.raise_for_status()
    res.encoding = None
    html = res.text
```



- 키워드를 입력하여 그에 해당하는 '톡톡' 게시판의 글들을 수집합니다. 최신순으로 게시글을 찾고, 게시물 각각에 대한 url을 우선 수집하며, 이 후 개별 url에 대한 본문 및 댓글 크롤링을 진행합니다.
- 긴 본문을 형태소 단위로 토큰화 진행해야 합니다. 또 해당 토큰들이 문장 내의 순서를 유지하도록 하여 추후 앞 뒤 관계에 따른 감성 분석을 진행 하고자 합니다.

커뮤니티 크롤링

02. 클리앙



setting - Windows 메모장

파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)

[CRAWLER]

KEYWORD = 쿠팡

[EXCEL]

RESULT_PATH = ./결과물.xlsx

	F		G	
12	댓글	종이8		2021-04-29 07:19:46
13	댓글	DavidMiles		2021-04-29 07:31:30
14	3 개시글	공녀어	1917	2021-04-28 17:50:48
15	댓글	예순어		2021-04-28 17:51:10
16	댓글	e58M		2021-04-28 17:52:12
17	댓글	anglelimit		2021-04-28 17:54:58
18	댓글	공녀어		2021-04-28 17:57:16
19	댓글	anglelimit		2021-04-28 18:21:06
20	댓글	비읍		2021-04-28 17:57:11
21	댓글	일산간지오빠		2021-04-28 17:58:18
22	댓글	중추는구미호		2021-04-28 18:06:13
23	댓글	헛지로		2021-04-28 18:15:38
24	댓글	몽골7		2021-04-28 18:52:06
25	댓글	mi200		2021-04-28 19:22:18
26	4 개시글	CHCGV	4529	2021-04-28 17:20:59

링크

내용 / 댓글

대단하세요 ㄷㄷ

https://www.clien.net/service/board/park/16102646?combine=true&q=%EC%8F%A0%ED%8C%A1&p=0&sort=rec

사람 끌어 넣는거로 통으로 그나마 제일 대

송백질이1님 부모님이 지금 저걸 사다달라고

사람끌어넣는다고 새벽배송한다고 하지만 1

송노불이타님 그런연도 있겠군요.

4월초 발할때 전기장판 고장나서 급하게 .

송희영지향님 당일 배송 완료는 저도 처음이

물류센터가 가까우면 그렇게 출매가 있어요.

USB A 젠더

https://www.clien.net/service/board/cm_mac/16102228?combine=true&q=%EC%8F%A0%ED%8C%A1&p=0&sort=

USB 커넥터 켜는 단순히 커넥터 모양만 비

저도 항상 외장하드를 연결해 두고 사용해서

답 주신 두 분 감사 드립니다.게이볼 여러개

https://www.clien.net/service/board/cm_iphonien/16162004?combine=true&q=%EC%8F%A0%ED%8C%A1&p=0&

오우 부럽습니다..

저도 현재 SE2 쓰고 13 기다리고있었는데 나

말씀을 살필 할 수 있나보네요?

@anglelimit님 밀레그럼 뽀뽀알림봇을 사용

송공녀어님 오 저도 해봐야겠어요감사합니

이 가격이면 고교./samsung family out

27	댓글	비읍		2021-04-28 17:57:11
28	댓글	일산간지오빠		2021-04-28 17:58:18
29	댓글	중추는구미호		2021-04-28 18:06:13
30	댓글	헛지로		2021-04-28 18:15:38
31	댓글	몽골7		2021-04-28 18:52:06
32	댓글	mi200		2021-04-28 19:22:18
33	4 개시글	CHCGV	4529	2021-04-28 17:20:59

링크

내용 / 댓글

쿠팡 아이폰 11 프로 일부모델 카드 50% 즉시할인하세요.

https://www.clien.net/service/board/park/16101942?combine=true&q=%EC%8F%A0%ED%8C%A1&p=0&sort=rec

이 가격이면 무조건 가여요

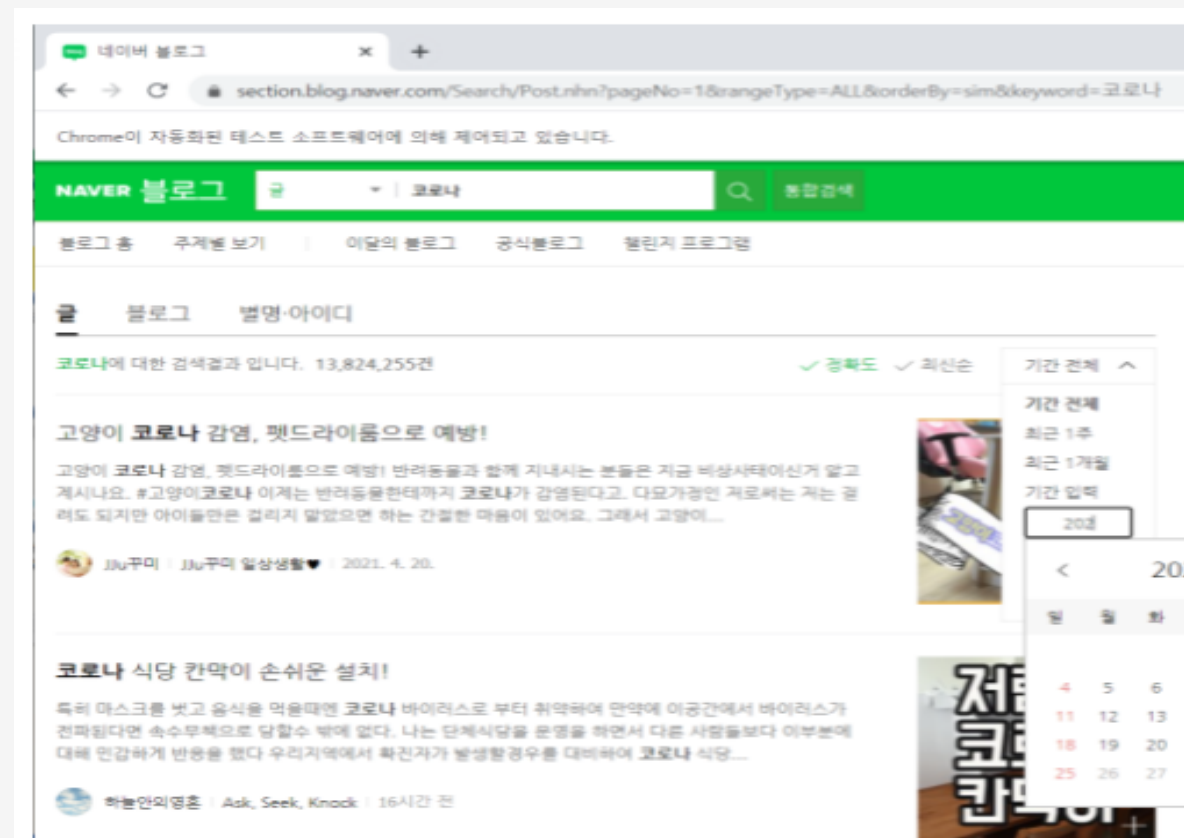
- 미리 txt 파일에 저장한 키워드에 대한 크롤링이 가능한 상태입니다.
- 프로그램 내에서 키워드 입력에 대한 크롤링으로 수정하고 앞서 회의를 통해 결정한 전처리 양식에 맞추어 크롤링한 데이터에 대한 정리가 필요합니다.

커뮤니티 외 크롤링 Test

03. 네이버 블로그 및 트위터

+

추후 커뮤니티 외 데이터가 많은 네이버 블로그 혹은 트위터(SNS)에서 감성분석에 적합한 데이터를 많이 얻을 수 있는지 test를 해 보았습니다.



- 데이터의 양이 매우 많고 또 최신 데이터가 많다는 장점이 있지만, 광고의 양도 많다는 단점이 있습니다.
=> 데이터 처리를 잘 한다면 활용가능한 자료가 될 것 입니다.



			username	tweet	retweets_count
88	2021-04-13	08:26:29	mbcnews	지난해에는 코로나19로 행사가 취소됐지만 올해는 예정대로 진행돼 수백 명의 참가자가...	10
89	2021-04-13	08:26:04	bissori0613	@CtWwBTS1 바다야🌊 거기 사는 모두에게 우리가 미안해🙏 코로나도 중요하지만...	0
90	2021-04-13	08:25:41	didchdkfl	코로나19에는 비대면 랜선라이프~ "달쌈광장"과 함께 핵 사이다 달쌈 맛집 gog...	0
91	2021-04-13	08:25:19	jason326naverc2	코로나	0
92	2021-04-13	08:24:54	kt114	논산여자중학교 제24기 운영위원회가 교목인 소나무를 기념식수 하며 출발 했네요. 최병...	0

- 자주 검색되는 단어의 경우 사람들이 노출 수를 올리기 위해 일부러 해시태그에 포함하기 때문에 관련 없는 트윗까지 크롤링 된다는 단점이 있습니다.
=> 주어지는 본문 글이 짧다는 단점이 있어 비교적 활용성이 떨어진다고 판단했습니다.

차주에는

1. 선정한 커뮤니티 4곳에 대한 크롤링부 (전처리를 포함한) 심화
 2. 크롤링 결과에 따라 어떤 감성어를 추가할 것인지 분석
 3. 감성어 사전 및 감성 분석 함수에 대한 학습 및 고민
- 을 진행하고자 합니다.



감사합니다 :)
빅데이터미네이터